

Minimax Lower Bounds for Noisy Matrix Completion Under Sparse Factor Models

Abhinav V. Sambasivan and Jarvis D. Haupt

Abstract

This paper examines fundamental error characteristics for a general class of matrix completion problems, where matrix of interest is a product of two a priori unknown matrices, one of which is sparse, and the observations are noisy. Our main contributions come in the form of minimax lower bounds for the expected per-element squared error for these problems under several noise/corruption models; specifically, we analyze scenarios where the corruptions are characterized by additive Gaussian noise or additive heavier-tailed (Laplace) noise, Poisson-distributed observations, and highly-quantized (e.g., one-bit) observations. Our results establish that the error bounds derived in (Soni et al., 2014) for *complexity-regularized maximum likelihood* estimators achieve, up to multiplicative constant and logarithmic factors, the minimax error rates in each of these noise scenarios, provided the sparse factor exhibits linear sparsity.

Index Terms

Matrix completion, dictionary learning, minimax lower bounds

I. INTRODUCTION

The matrix completion problem involves imputing the missing values of a matrix from an incomplete, and possibly noisy sampling of its entries. In general, without making any assumption about the entries of the matrix, the matrix completion problem is ill-posed and it is impossible to recover the matrix uniquely. However, if the matrix to be recovered has some intrinsic structure (e.g., low rank structure), it is possible to design algorithms which exactly estimate the missing entries. Indeed, the performance of convex methods for low-rank matrix completion problems have been extensively studied in noiseless settings [1]–[5], in noisy settings where the observations are affected by additive noise [6]–[12], and in settings where the observations are non-linear (e.g., highly-quantized or Poisson distributed observation) functions of the underlying matrix entry (see, [13]–[15]). Recent works which explore robust recovery of low-rank matrices under malicious corruptions which are sparse include [16]–[19].

A notable advantage of using low-rank models is that the estimation strategies involved in completing such matrices can be cast into efficient convex methods which are well-understood and suitable to analyses. The fundamental estimation error characteristics for more general completion problems, for example, those employing

Submitted September 28, 2015. The authors are with the Department of Electrical and Computer Engineering at the University of Minnesota – Twin Cities. Tel/fax: (612) 625-3300 / (612) 625-4583. Emails: {samba014, jdhaupt}@umn.edu. The authors graciously acknowledge support from the DARPA Young Faculty Award, Grant No. N66001-14-1-4047.

general bilinear factor models, have not (to our knowledge) been fully characterized. In this work we provide several new results in this direction. Our focus here is on matrix completion problems under *sparse factor model* assumptions, where the matrix to be estimated is well-approximated by a product of two matrices, one of which is sparse. Such models have been motivated by a variety of applications in dictionary learning, subspace clustering, image demosaicing, and various machine learning problems (see, e.g. the discussion in [20]). Here, we investigate fundamental lower bounds on the achievable estimation error for these problems in several specific noise scenarios – additive Gaussian noise, additive heavier-tailed (Laplace) noise, Poisson-distributed observations, and highly-quantized (e.g., one-bit) observations. Our analyses compliment the upper bounds provided recently in [20] for complexity-penalized maximum likelihood estimation methods, and establish that the error rates obtained in [20] are nearly minimax optimal under an (arguably, natural) assumption, that the sparse factor exhibit *linear* sparsity (so that its number of non zeros is a constant fraction of its size).

The remainder of the paper is organized as follows. We begin with a brief overview of the various preliminaries and notations and a formal definition of the matrix completion problem considered here in Section II. Our main results are stated in Section III; there, we establish minimax lower bounds for the recovery of a matrix \mathbf{X} which admits a sparse factorization under various noise models, and also briefly discuss the implications of these bounds and compare them with existing works. In Section IV we conclude with a concise discussion of possible extensions and potential future directions. The proofs of our main results are provided in the Appendix.

A. Notations and Preliminaries

We provide a brief summary of the notations used here and revisit a few key concepts before delving into our main results.

We let $a \vee b = \max\{a, b\}$ and $a \wedge b = \min\{a, b\}$. For any $n \in \mathbb{N}$, $[n]$ denotes the set of integers $\{1, \dots, n\}$. For a matrix \mathbf{M} , we use the following notation: $\|\mathbf{M}\|_0$ denotes the number of non-zero elements in \mathbf{M} , $\|\mathbf{M}\|_\infty = \max_{i,j} |M_{i,j}|$ denotes the entry-wise maximum (absolute) entry of \mathbf{M} , $\|\mathbf{M}\|_F = \sqrt{\sum_{i,j} M_{i,j}^2}$ denotes the Frobenius norm. We use the standard asymptotic computational complexity (Big O and Big Omega) notations to suppress leading constants in our results for, clarity of exposition.

We also briefly recall an important information-theoretic quantity, the Kullback-Leibler divergence (or KL divergence). When $x(z)$ and $y(z)$ denote the pdf (or pmf) of a real valued random variable Z , the KL divergence of y from x is denoted by $K(\mathbb{P}_x, \mathbb{P}_y)$ and given by

$$\begin{aligned} K(\mathbb{P}_x, \mathbb{P}_y) &= \mathbb{E}_{Z \sim x(Z)} \left[\log \frac{x(Z)}{y(Z)} \right] \\ &\triangleq \mathbb{E}_x \left[\log \frac{x(Z)}{y(Z)} \right], \end{aligned}$$

provided $x(z) = 0$ whenever $y(z) = 0$, and ∞ otherwise. The logarithm is taken to be the natural log. It is worth noting that the KL divergence is not necessarily commutative and, $K(\mathbb{P}_x, \mathbb{P}_y) \geq 0$ with $K(\mathbb{P}_x, \mathbb{P}_y) = 0$ when $x(Z) = y(Z)$. In a sense, the KL divergence quantifies how “far” apart two distributions are.

II. PROBLEM STATEMENT

A. Observation Model

We consider the problem of estimating the entries of an unknown matrix $\mathbf{X}^* \in \mathbb{R}^{n_1 \times n_2}$ which admits a factorization of the form,

$$\mathbf{X}^* = \mathbf{D}^* \mathbf{A}^*, \quad (1)$$

where for some integer $1 \leq r \leq (n_1 \wedge n_2)$, $\mathbf{D}^* \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{A}^* \in \mathbb{R}^{r \times n_2}$ are *a priori* unknown factors. Additionally, our focus in this paper will be restricted to the cases where the matrix \mathbf{A}^* is k -sparse (having no more than $0 < k \leq rn_2$ nonzero elements). We further assume that the elements of \mathbf{D}^* and \mathbf{A}^* are bounded, i.e.

$$\|\mathbf{D}^*\|_\infty \leq 1 \text{ and } \|\mathbf{A}^*\|_\infty \leq A_{\max} \quad (2)$$

for some constant $A_{\max} > 0$. A direct implication of (2) is that the elements of \mathbf{X}^* are also bounded (with $\|\mathbf{X}^*\|_\infty \leq X_{\max} \leq rA_{\max}$). While bounds on the amplitudes of the elements of the matrix to be estimated often arise naturally in practice, the assumption that the entries of the factors are bounded fixes some of the scaling ambiguities associated with the bilinear model.

Instead of observing all the elements of the matrix \mathbf{X}^* directly, we assume here that we make noisy observations of \mathbf{X}^* at a known *subset* of locations. In what follows, we will model the observations $Y_{i,j}$ as i.i.d draws from a probability distribution (or mass) function parameterized by the true underlying matrix entry $X_{i,j}^*$. We denote by $\mathcal{S} \subseteq [n_1] \times [n_2]$ the set of locations at which observations are collected, and assume that these points are sampled uniformly with $\mathbb{E}[|\mathcal{S}|] = m$ (which denotes the nominal number of measurements) for some integer m satisfying $1 \leq m \leq n_1 n_2$. Specifically, for $\gamma = m/(n_1 n_2)$, we suppose \mathcal{S} is generated according to an independent Bernoulli(γ) model, so that each $(i, j) \in [n_1] \times [n_2]$ is included in \mathcal{S} independently with probability γ . Thus, given \mathcal{S} , we model the collection of $|\mathcal{S}|$ measurements of \mathbf{X}^* in terms of the collection $\{Y_{i,j}\}_{(i,j) \in \mathcal{S}} \triangleq \mathbf{Y}_{\mathcal{S}}$ of conditionally (on \mathcal{S}) independent random quantities. The joint pdf (or pmf) of the observations can be formally written as

$$p_{\mathbf{X}_{\mathcal{S}}^*}(\mathbf{Y}_{\mathcal{S}}) \triangleq \prod_{(i,j) \in \mathcal{S}} p_{X_{i,j}^*}(Y_{i,j}) \triangleq \mathbb{P}_{\mathbf{X}^*}, \quad (3)$$

where $p_{X_{i,j}^*}(Y_{i,j})$ denotes the corresponding scalar pdf (or pmf), and we use the shorthand $\mathbf{X}_{\mathcal{S}}^*$ to denote the collection of elements of \mathbf{X}^* indexed by $(i, j) \in \mathcal{S}$. Given \mathcal{S} and the corresponding noisy observations $\mathbf{Y}_{\mathcal{S}}$ of \mathbf{X}^* distributed according to (3), the matrix completion problem aims at estimating \mathbf{X}^* under the assumption that it admits a sparse factorization as in (1).

B. The minimax risk

In this paper, we examine the fundamental limits of estimating the elements of a matrix which follows the model (1) and observations as described above, using any possible estimator (irrespective of its computational tractability).

The accuracy of an estimator $\hat{\mathbf{X}}$ in estimating the entries of the true matrix \mathbf{X}^* is measured in terms of its risk $\mathcal{R}_{\hat{\mathbf{X}}}$ which we define to be the normalized (per-element) Frobenius error,

$$\mathcal{R}_{\hat{\mathbf{X}}} = \frac{\mathbb{E}_{\mathbf{Y}_S} [\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2]}{n_1 n_2}. \quad (4)$$

Here, our notation is meant to denote that the expectation is taken with respect to all of the random quantities (i.e., the joint distribution of \mathcal{S} and \mathbf{Y}_S).

Let us now consider a class of matrices parameterized by the inner dimension r , sparsity factor k and upper bound on the amplitude of elements of \mathbf{A} , A_{\max} , where each element in the class obeys the factor model (1) and the assumptions in (2). Formally, we set

$$\mathcal{X}(r, k, A_{\max}) \triangleq \{\mathbf{X} = \mathbf{D}\mathbf{A} \in \mathbb{R}^{n_1 \times n_2} : \mathbf{D} \in \mathbb{R}^{n_1 \times r}, \|\mathbf{D}\|_{\infty} \leq 1 \text{ and } \mathbf{A} \in \mathbb{R}^{r \times n_2}, \|\mathbf{A}\|_0 \leq k, \|\mathbf{A}\|_{\infty} \leq A_{\max}\}. \quad (5)$$

The worst-case performance of an estimator $\hat{\mathbf{X}}$ over the class $\mathcal{X}(r, k, A_{\max})$, under the Frobenius error metric defined in (4), is given by its maximum risk,

$$\tilde{\mathcal{R}}_{\hat{\mathbf{X}}} \triangleq \sup_{\mathbf{X}^* \in \mathcal{X}(r, k, A_{\max})} \mathcal{R}_{\hat{\mathbf{X}}}.$$

The estimator having the smallest maximum risk among all possible estimators and is said to achieve the minimax risk, which is a characteristic of the estimation problem itself. For the problem of matrix completion under the sparse factor model described in Section II-A, the minimax risk is expressed as

$$\begin{aligned} \mathcal{R}_{\mathcal{X}(r, k, A_{\max})}^* &\triangleq \inf_{\hat{\mathbf{X}}} \tilde{\mathcal{R}}_{\hat{\mathbf{X}}} \\ &= \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}(r, k, A_{\max})} \mathcal{R}_{\hat{\mathbf{X}}} \\ &= \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}(r, k, A_{\max})} \frac{\mathbb{E}_{\mathbf{Y}_S} [\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2]}{n_1 n_2}. \end{aligned} \quad (6)$$

As we see, the minimax risk depends on the choice of the model class parameters r , k and A_{\max} . It is worth noting that inherent in the formulation of the minimax risk are the noise model and the nominal number of observations ($m = \mathbb{E}[|\mathcal{S}|]$) made. For the sake of brevity, we shall not make all such dependencies explicit.

In general it is complicated to obtain closed form solutions for (6). Here, we will adopt a common approach employed for such problems, and seek to obtain lower bounds on the minimax risk, $\mathcal{R}_{\mathcal{X}(r, k, A_{\max})}^*$ using tools from [21]. Our analytical approach is inspired also by the approach in [22], which considered the problem of estimating low-rank matrices corrupted by sparse outliers.

III. MAIN RESULTS AND IMPLICATIONS

In this section we establish lower bounds on the minimax risk for the problem settings defined in Section II for different noise models and various regimes of sparsity. The two sparsity regimes examined include

- $k < n_2$: There is (on an average), less than one non-zero element per column of \mathbf{A}^*
- $k \geq n_2$: There is (on an average), more than one non-zero element per column of \mathbf{A}^*

A. Additive Gaussian Noise

Let us consider a setting where the observations are corrupted by i.i.d zero-mean additive Gaussian noise with known variance. We have the following result; its proof appears in Appendix A

Theorem III.1. *For observations made as per the model, $Y_{i,j} = X_{i,j}^* + \xi_{i,j}$, suppose that the variables $\xi_{i,j}$ are i.i.d Gaussian $\mathcal{N}(0, \sigma^2)$, $\sigma > 0$, $\forall (i, j) \in \mathcal{S}$. Then, there exists an absolute constant $c > 0$ such that for all $n_1, n_2 \geq 2$ and $1 \leq r \leq (n_1 \wedge n_2)$, the minimax risk for sparse factor matrix completion over the model class $\mathcal{X}(r, k, A_{\max})$ obeys*

$$\mathcal{R}_{\mathcal{X}(r, k, A_{\max})}^* \geq c (\sigma \wedge A_{\max})^2 \left\{ \left(\frac{n_1 r}{m} \right) \Delta(k, n_2) + \left(\frac{k}{m} \right) \left(\frac{k}{r n_2} \right) \right\}, \quad (7)$$

where the function $\Delta(k, n_2)$ is given by

$$\Delta(k, n_2) = \begin{cases} (k/n_2) & \text{for } k < n_2 \\ 1 & \text{for } k \geq n_2 \end{cases}. \quad (8)$$

Remark III.1. *If instead of i.i.d Gaussian noise, we have that $\xi_{i,j}$ are just independent zero-mean additive Gaussian random variables with variances $\sigma_{i,j}^2 \geq \sigma_{\min}^2 \forall (i, j) \in \mathcal{S}$, the result in (7) is still valid with the σ replaced by σ_{\min} . This stems from the fact that the KL divergence between the distributions in equations (39) and (41) can be easily upper bounded by smallest of value of variance amongst all the noise entries.*

Let us now analyze the result of this theorem more closely and see how the estimation risk varies as a function of the number of measurements obtained, as well as the dimension and sparsity parameters of the matrix to be estimated. The minimax risk as in equation (7) consists of two main terms

- The $\left(\frac{n_1 r}{m} \right) \Delta(k, n_2)$ term may be interpreted as the error associated with estimating the non sparse factor \mathbf{D}^* . Here the quantity $n_1 r$ which gives the number of elements in \mathbf{D}^* , can be viewed as the number of *degrees of freedom* contributed by the non-sparse factor in the matrix to be estimated. It can be seen that error associated with the non-sparse factor follows the parametric rate $(n_1 r/m)$ when $k \geq n_2$, i.e. \mathbf{A}^* (on an average) has more than one non-zero element per column. Qualitatively, this implies that all the degrees of freedom offered by \mathbf{D}^* manifest in the estimation of the overall matrix \mathbf{X}^* provided there are enough non-zero elements (at least one non-zero per column) in \mathbf{A}^* . If there are (on an average) less than one non-zero element per column in the sparse factor, a few rows of \mathbf{D}^* vanish due to the presence of zero columns in \mathbf{A}^* and hence all the degrees of freedom in \mathbf{D}^* are not carried over to \mathbf{X}^* . This makes the overall problem easier and reduces the minimax risk by a factor of (k/n_2) .
- The $\left(\frac{k}{m} \right) \left(\frac{k}{r n_2} \right)$ term may be interpreted as the error associated with estimating the sparse factor \mathbf{A}^* . As in the previous case, the quantity k which indicates the maximum number of non-zero elements that could be present in the sparse matrix \mathbf{A}^* can be interpreted as the number of *degrees of freedom* contributed by the sparse factor in the matrix to be estimated. The error associated with the sparse factor, follows the parametric rate (k/m) scaled by the relative sparsity $(k/r n_2)$, (ratio of maximum number of non-zero elements to the

total number of elements) of \mathbf{A}^* . Qualitatively speaking, the relative sparsity term determines by how much do the degrees of freedom rendered by \mathbf{A}^* is reflected in the estimation of the overall product matrix \mathbf{X}^* .

In the following remark we shall consider a specific instance of the result in (7) when we make certain assumptions about the sparsity parameter k and then discuss its relation to already existing work in matrix completion.

Remark III.2. *When the factor \mathbf{A}^* has linear sparsity i.e. if $\exists c' \in (0, 1)$ such that $k = c'rn_2$, then the minimax risk obtained in (7) obeys the lower bound*

$$\mathcal{R}_{\mathcal{X}(r,k,A_{\max})}^* = \Omega\left((\sigma \wedge A_{\max})^2 \left(\frac{n_1 r + k}{m}\right)\right), \quad (9)$$

where the $\Omega(\cdot)$ notation suppresses leading constants and illustrates the dependence of the minimax risk in terms of the key problem parameters.

It is worth noting that the minimax risk in the linearly sparse regime relates directly to the work in [20], which gives upper bounds for matrix completion problems under similar sparse factor models. The normalized (per-element) Frobenius error for the sparsity-penalized maximum likelihood estimator under a Gaussian noise model presented in [20] satisfies

$$\frac{\mathbb{E}_{\mathbf{Y}_S} \left[\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \right]}{n_1 n_2} = \mathcal{O} \left((\sigma \wedge A_{\max})^2 \left(\frac{n_1 r + k}{m} \right) \log(n_1 \vee n_2) \right). \quad (10)$$

A comparison of (10) to our results in Equations (7) and (9) imply that the rate attained by the estimator presented in [20] is minimax optimal up to a logarithmic factor in the linearly sparse regime.

Another direct point of comparison to our result here is the low rank matrix completion problem with entry-wise observations considered in [8]. In particular if we adopt the lower bounds obtained in Theorem 6 of their work to our settings, we observe that the risk involved in estimating rank- r matrices which are sampled uniformly at random follows

$$\begin{aligned} \frac{\mathbb{E}_{\mathbf{Y}_S} \left[\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \right]}{n_1 n_2} &= \Omega \left((\sigma \wedge X_{\max})^2 \left(\frac{(n_1 \vee n_2)r}{m} \right) \right) \\ &= \Omega \left((\sigma \wedge X_{\max})^2 \left(\frac{(n_1 + n_2)r}{m} \right) \right), \end{aligned} \quad (11)$$

where the last inequality follows from the fact that $n_1 \vee n_2 \geq (n_1 + n_2)/2$. If we consider non-sparse factor models (where $k = rn_2$), it can be seen that the product $\mathbf{X}^* = \mathbf{D}^* \mathbf{A}^*$ is low-rank with $\text{rank}(\mathbf{X}^*) \leq r$ and our problem reduces to the one considered in [8]. Under the conditions described above, the lower bound given by us in (7) coincides (up to a constant scaling factor) with (11), provided $X_{\max} = \Omega(A_{\max})$. However the introduction of sparsity brings additional structure which can be exploited in estimating the entries of \mathbf{X}^* thus decreasing the risk involved.

B. Additive Laplace Noise

Suppose the observations \mathbf{Y}_S are corrupted with heavier tailed noises like i.i.d zero-mean additive Laplace noise with parameter $\tau > 0$. The following theorem gives a lower bound on the minimax risk in such settings; a sketch

of the proof is given in Appendix B.

Theorem III.2. *For observations made as per the model $Y_{i,j} = X_{i,j}^* + \xi_{i,j}$, suppose that the variables $\xi_{i,j}$ are i.i.d Laplace(0, τ), $\tau > 0$, $\forall (i, j) \in \mathcal{S}$. Then, there exists an absolute constant $c > 0$ such that for all $n_1, n_2 \geq 2$ and $1 \leq r \leq (n_1 \wedge n_2)$, the minimax risk for sparse factor matrix completion over the model class $\mathcal{X}(r, k, A_{\max})$ obeys*

$$\mathcal{R}_{\mathcal{X}(r, k, A_{\max})}^* \geq c (\tau^{-1} \wedge A_{\max})^2 \left\{ \left(\frac{n_1 r}{m} \right) \Delta(k, n_2) + \left(\frac{k}{m} \right) \left(\frac{k}{r n_2} \right) \right\}, \quad (12)$$

where the function $\Delta(k, n_2)$ depends on the sparsity regime and is defined in (8).

When we compare the lower bounds obtained under this noise model to the results of the previous case it can be readily seen that the overall error rates achieved is similar in both cases. Since we have the variance of Laplace(τ) random variable to be $(2/\tau^2)$, the leading term $(\tau^{-1} \wedge A_{\max})^2$ here is analogous to the $(\sigma \wedge A_{\max})^2$ factor which appears in the error bound for Gaussian noise. For linearly sparse regimes discussed in Remark III.2, the minimax risk can be easily shown to attain parametric rates similar to (9). Using (12), we can observe that the complexity penalized maximum likelihood estimator described in [20] is minimax optimal up to a constant times a logarithmic factor, $\tau X_{\max} \log(n_1 \vee n_2)$.

C. Poisson-distributed Observations

Let us now consider a scenario where the data maybe observed as discrete ‘counts’ (which is common in imaging applications e.g., number of photons hitting the receiver per unit time). A popular model for such settings is the Poisson-distributed observation model where all the entries of the matrix \mathbf{X}^* to be estimated are positive and our observation $Y_{i,j}$ at each location $(i, j) \in \mathcal{S}$ is an independent Poisson random variable with a rate parameter $X_{i,j}^*$. The problem of matrix completion now involves the task of Poisson denoising. In this case, we get the following result; its proof appears in Appendix C.

Theorem III.3. *Suppose that the entries of the matrix \mathbf{X}^* satisfy $\min_{(i,j) \in [n_1] \times [n_2]} X_{i,j} \geq X_{\min}$ for some constant $0 < X_{\min} < A_{\max}$ and the observations $Y_{i,j}$ are independent Poisson distributed random variable with rates $X_{i,j}^* \forall (i, j) \in \mathcal{S}$. Then, there exists an absolute constant $c > 0$ such that for all $n_1, n_2 \geq 2$ and $1 \leq r \leq (n_1 \wedge n_2)$, the minimax risk for sparse factor matrix completion over the model class $\mathcal{X}'(r, k, A_{\max})$ which is a subset of $\mathcal{X}(r, k, A_{\max})$ comprised of matrices with positive entries, obeys*

$$\mathcal{R}_{\mathcal{X}'(r, k, A_{\max})}^* \geq c \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right]^2 \left\{ \left(\frac{n_1 r}{m} \right) + \left(\frac{k - n_2}{m} \right) \left(\frac{k - n_2}{r n_2} \right) \right\}. \quad (13)$$

As in the previous case, our analysis rests on establishing quadratic upper bounds on the KL divergence to obtain parametric error rates for the minimax risk; a similar approach was used in [23], which describes performance bounds on compressive sensing sparse signal estimation task under a Poisson noise model, and in [24]. Recall that the lower bounds for each of the preceding cases exhibited a leading factor which was essentially the minimum of the noise variance and A_{\max}^2 . Unlike those cases, for a Poisson observation model, the noise variance equals the

rate parameter and hence depends on the true underlying matrix entry. So, for $X_{\min} \ll 1$, we might interpret the factor $(X_{\min} + \sqrt{X_{\min}})^2 \approx X_{\min}$, which is the minimum variance of all the Poisson distributed observations and hence is analogous to the results presented for the Gaussian and Laplace noise models.

The dependence of the minimax risk on the nominal number of observations (m), matrix dimensions (n_1, n_2, r), and sparsity factor k , is encapsulated in the two terms, $\left(\frac{n_1 r}{m}\right)$ and $\left(\frac{k-n_2}{m}\right)\left(\frac{k-n_2}{r n_2}\right)$. The first term which corresponds to the error associated with the dictionary term \mathbf{D}^* is exactly the same as in the previous noise models. However we can see that the term associated with the sparse factor \mathbf{A}^* is a bit different from the other models discussed. In a Poisson-distributed observation model, we have that the entries of the true underlying matrix to be estimated are positive (which also serves as the Poisson rate parameter to the observations $Y_{i,j}$). A necessary implication of this is that the sparse factor \mathbf{A}^* should contain no zero-valued columns, or every columns should have at least one non-zero entry (and hence we have $k \geq n_2$). This reduces the number of *degrees of freedom* (as described in Section III-A) in the sparse factor from k to $k - n_2$, thus reducing the over all minimax risk.

It is worth further commenting on the relevance of this result to the work in [20], which establishes error bounds for Poisson denoising problems with sparse factor models. Casting the results of that work to our settings, we see that the normalized (per element) error of the complexity penalized maximum likelihood estimator obeys

$$\frac{\mathbb{E}_{\mathbf{Y}_S} \left[\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2 \right]}{n_1 n_2} = \mathcal{O} \left(\left(X_{\max} + \frac{X_{\max}}{X_{\min}} \cdot X_{\max}^2 \right) \left(\frac{n_1 r + k}{m} \right) \log(n_1 \vee n_2) \right), \quad (14)$$

where the X_{\max} ($\leq r A_{\max}$) term is the upper bound on the entries of the matrix to be estimated. Comparing (14) with the lower bound established in (13), we can see that estimator described in [20] is minimax optimal w.r.t to the matrix dimension parameters up to a logarithmic factor (neglecting the leading constants) in the linearly sparse regime.

We comment a bit on our assumption that the elements of the true underlying matrix \mathbf{X}^* , be greater than or equal to some $X_{\min} > 0$. Here, this parameter shows up in the leading term as $(X_{\min} + \sqrt{X_{\min}})$, which suggests that the minimax risk vanishes as the rate of the Poisson processes tends to 0. This implication is in agreement with the Cramér-Rao lower bounds which states that the error associated in estimating a $\text{Poisson}(\theta)$ random variable using n iid observations decays at the rate θ/n (which is achieved by a sample average estimator). Thus our notion that the denoising problem getting *easier* as the rate parameter decreases is intuitive and consistent with classical analyses. On this note, we briefly mention recent efforts which do not make assumptions on the minimum rate of the underlying Poisson processes; for matrix estimation tasks as here [15], and for sparse vector estimation from Poisson-distributed compressive observations [25].

D. One-bit Observation Model

We consider here a scenario where the observations are quantized to a single bit i.e. the observations $Y_{i,j}$ can take only binary values (either 0 or 1). Quantized observation models arise in many collaborative filtering applications where the user ratings are quantized to fixed levels, in quantum physics, communication networks, etc. (see, e.g. discussions in [13], [26]).

For a given sampling set \mathcal{S} , we consider the observations $\mathbf{Y}_{\mathcal{S}}$ to be conditionally (on \mathcal{S}) independent random quantities defined by

$$Y_{i,j} = \mathbf{1}_{\{Z_{i,j} \geq 0\}}, \quad (i,j) \in \mathcal{S}, \quad (15)$$

where

$$Z_{i,j} = X_{i,j}^* - W_{i,j},$$

the $\{W_{i,j}\}_{(i,j) \in \mathcal{S}}$ are i.i.d continuous zero-mean scalar noises having (bounded) probability density function $f(w)$ and cumulative density function $F(w)$ for $w \in \mathbb{R}$, and $\mathbf{1}_{\{\mathcal{A}\}}$ is the indicator function which takes the value 1 when the event \mathcal{A} occurs (or is true) and zero otherwise. Thus our observations are quantized corrupted versions of the true underlying matrix entries. Note that the independence of $W_{i,j}$ implies that the elements $Y_{i,j}$ are also independent. Given this model, it can be easily seen that each $Y_{i,j} \forall (i,j) \in \mathcal{S}$ is a Bernoulli random variable whose parameter is a function of the true parameter $X_{i,j}^*$, and the cumulative density function $F(\cdot)$. In particular, for any $(i,j) \in \mathcal{S}$, we have $\Pr(Y_{i,j} = 1) = \Pr(W_{i,j} \leq X_{i,j}^*) = F(X_{i,j}^*)$. Hence the joint pmf of the observations $\mathbf{Y}_{\mathcal{S}} \in \{0,1\}^{|\mathcal{S}|}$ (conditioned on the underlying matrix entries) can be written as,

$$p_{\mathbf{X}_{\mathcal{S}}}^*(\mathbf{Y}_{\mathcal{S}}) = \prod_{(i,j) \in \mathcal{S}} [F(X_{i,j}^*)]^{Y_{i,j}} [1 - F(X_{i,j}^*)]^{1-Y_{i,j}}. \quad (16)$$

We will further assume that $F(rA_{\max}) < 1$ and $F(-rA_{\max}) > 0$, which will allow us to avoid some pathological scenarios in our results. In such settings, the following theorem gives a lower bound on the minimax risk; a sketch of the proof is given in Appendix D.

Theorem III.4. *Suppose that the observations $Y_{i,j}$ are obtained as described in (15), where $W_{i,j}$ are iid continuous zero-mean scalar random variables $\forall (i,j) \in \mathcal{S}$ having probability density function $f(w)$, and cumulative density function $F(w)$ for $w \in \mathbb{R}$. Define*

$$c_{F,rA_{\max}} \triangleq \left(\sup_{|t| \leq rA_{\max}} \frac{1}{F(t)(1-F(t))} \right)^{1/2} \left(\sup_{|t| \leq rA_{\max}} f^2(t) \right)^{1/2}. \quad (17)$$

Then, there exists an absolute constant $c > 0$ such that for all $n_1, n_2 \geq 2$ and $1 \leq r \leq (n_1 \wedge n_2)$, the minimax risk for sparse factor matrix completion over the model class $\mathcal{X}(r, k, A_{\max})$ obeys

$$\mathcal{R}_{\mathcal{X}(r,k,A_{\max})}^* \geq c \left[c_{F,rA_{\max}}^{-1} \wedge A_{\max} \right]^2 \left\{ \left(\frac{n_1 r}{m} \right) \Delta(k, n_2) + \left(\frac{k}{m} \right) \left(\frac{k}{rn_2} \right) \right\}, \quad (18)$$

where the function $\Delta(k, n_2)$ depends on the sparsity regime and is defined in Equation (8).

It worth commenting on the relevance of our result (in the linear sparsity regime) to the upper bounds established in [20], for the matrix completion problem under similar settings. The normalized (per element) error of the complexity penalized maximum likelihood estimator described in [20] obeys

$$\frac{\mathbb{E}_{\mathbf{Y}_{\mathcal{S}}} [\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2]}{n_1 n_2} = \mathcal{O} \left(\left(\frac{c_{F,rA_{\max}}^2}{c_{F,rA_{\max}}'} \right) \left(\frac{1}{c_{F,rA_{\max}}^2} + X_{\max}^2 \right) \left(\frac{n_1 r + k}{m} \right) \log(n_1 \vee n_2) \right), \quad (19)$$

where $X_{\max} (\geq 0)$ is the upper bound on the entries of the matrix to be estimated and $c'_{F, rA_{\max}}$ is defined as

$$c'_{F, rA_{\max}} \triangleq \inf_{|t| \leq rA_{\max}} \frac{f^2(t)}{F(t)(1-F(t))}. \quad (20)$$

Comparing (19) with the lower bound established in (18), we can see that estimator described in [20] is minimax optimal up to a logarithmic factor (in the linearly sparse regime) when the term $(c_{F, rA_{\max}}^2 / c'_{F, rA_{\max}})$ is bounded above by a constant. The lower bounds obtained for the one-bit observation model and the Gaussian case essentially exhibit the same dependence on the matrix dimensions (n_1, n_2) and r , sparsity (k) and the nominal number of measurements (m) , except for the leading term (which explicitly depends on the distribution of the noise variables $W_{i,j}$ for the one-bit case). Such a dependence in error rates between rate-constrained tasks and their Gaussian counterparts was observed in earlier works on rate-constrained parameter estimation [27], [28].

It is also interesting to compare our result with the lower bounds for the low-rank matrix completion problem considered in [13]. In that work, the authors establish that the risk involved in matrix completion over a (convex) set of max-norm and nuclear norm constrained matrices (with the decreasing noise pdf $f(t)$ for $t > 0$) obeys

$$\begin{aligned} \frac{\mathbb{E}_{\mathbf{Y}_S} [\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2]}{n_1 n_2} &= \Omega \left(X_{\max} \sqrt{\frac{1}{c'_{F, rA_{\max}}}} \sqrt{\frac{(n_1 \vee n_2)r}{m}} \right) \\ &= \Omega \left(X_{\max} \sqrt{\frac{1}{c'_{F, rA_{\max}}}} \sqrt{\frac{(n_1 + n_2)r}{m}} \right), \end{aligned} \quad (21)$$

where $c'_{F, rA_{\max}}$ is defined as in (20). As long as $c_{F, rA_{\max}}^2$ and $\sqrt{c'_{F, rA_{\max}}}$ are comparable, the leading terms of the two bounds are analogous to each other. In order to note the difference between this result and ours, we consider the case when \mathbf{A}^* is not sparse i.e., we set $k = rn_2$ in (18) so that the resulting matrix \mathbf{X}^* is low-rank (with $\text{rank}(\mathbf{X}) \leq r$). For such a setting, our error bound (18) scales in proportion to the ratio of the degrees of freedom $(n_1 + n_2)r$ and the nominal number of observations m , while the bound in [13] scales to the square root of that ratio.

A more recent work [29], proposed an estimator for the low-rank matrix completion on finite alphabets and establishes convergence rates faster than in [13]. On casting their results to our settings, the estimation error in [29] was shown to obey

$$\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} = \mathcal{O} \left(\left(\frac{c_{F, rA_{\max}}^2}{c'_{F, rA_{\max}}} \right)^2 \left(\frac{(n_1 + n_2)r}{m} \right) \log(n_1 + n_2) \right). \quad (22)$$

On comparing (22) with the our lower bounds (for the low-rank case, where $k = rn_2$), it is worth noting that their estimator achieves minimax optimal rates up to a logarithmic factor when the ratio $(c_{F, rA_{\max}}^2 / c'_{F, rA_{\max}})$ is bounded above by a constant.

IV. CONCLUSION

In this paper, we established minimax lower bounds for sparse factor matrix completion tasks, under several different noise/corruption models. It is interesting to note that, while our focus here was on several specific noise

models, the essential structure of our analysis could be easily extended to any model for which the scalar KL divergence may be upper bounded by a function quadratic in the parameter difference.

A unique aspect of our analysis is its applicability to matrices representable as a product of structured factors. While our focus here was specifically on models in which one factor is sparse, the approach we utilize here may be extended to other structured factor models (of which standard low-rank models are one particular case). A similar analysis to that utilized here could also be used to establish lower bounds on estimation of structured tensors, for example, those expressible in a Tucker decomposition with sparse core, and possibly structured factor matrices (see, e.g., [30] for a discussion of Tucker models). We defer investigations along these lines to a future effort.

APPENDIX

In order to prove Theorems III.1 to III.4 we use standard minimax analysis techniques from [22], namely the following theorem (whose proof is available in [21]),

Theorem A.1 (Adopted from Theorem 2.5 in [21]). *Assume that $M \geq 2$ and suppose that there exists a set with finite elements, $\mathcal{X} = \{\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_M\} \subset \mathcal{X}(r, k, A_{\max})$ such that*

- $d(\mathbf{X}_j, \mathbf{X}_k) \geq 2s, \quad \forall 0 \leq j < k \leq M$; where $d(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a semi-distance function, and
- $\frac{1}{M} \sum_{j=1}^M K(\mathbb{P}_{\mathbf{X}_j}, \mathbb{P}_{\mathbf{X}_0}) \leq \alpha \log M$ with $0 < \alpha < 1/8$.

Then

$$\begin{aligned} \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{X}(r, k, A_{\max})} \mathbb{P}_{\mathbf{X}}(d(\hat{\mathbf{X}}, \mathbf{X}) \geq s) &\geq \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{X}} \mathbb{P}_{\mathbf{X}}(d(\hat{\mathbf{X}}, \mathbf{X}) \geq s) \\ &\geq \frac{\sqrt{M}}{1 + \sqrt{M}} \left(1 - 2\alpha - \sqrt{\frac{2\alpha}{\log M}} \right) > 0. \end{aligned} \quad (23)$$

Here the first inequality arises from the fact that the supremum over a class of matrices \mathcal{X} is upper bounded by that of a larger class $\mathcal{X}(r, k, A_{\max})$ (or in other words, estimating the matrix over an uncountably infinite class is at least as difficult as solving the problem over any finite subclass). We thus reduce the problem of matrix completion over an uncountably infinite set $\mathcal{X}(r, k, A_{\max})$, to a carefully chosen finite collection of matrices $\mathcal{X} \subset \mathcal{X}(r, k, A_{\max})$ and lower bound the latter which then gives a valid bound for the overall problem.

A. Proof of Theorem III.1

Let us define a class of matrices $\mathcal{X} \subset \mathbb{R}^{n_1 \times n_2}$ as

$$\mathcal{X} \triangleq \{\mathbf{X} = \mathbf{D}\mathbf{A} : \mathbf{D} \in \mathcal{D}, \mathbf{A} \in \mathcal{A}\}, \quad (24)$$

where the factor classes $\mathcal{D} \subset \mathbb{R}^{n_1 \times r}$ and $\mathcal{A} \subset \mathbb{R}^{r \times n_2}$ are constructed as follows for $\gamma_d, \gamma_a \leq 1$ (which we shall qualify later)

$$\mathcal{D} = \left\{ \mathbf{D} = (d_{ij}) \in \mathbb{R}^{n_1 \times r} : d_{ij} \in \left\{ 0, 1, \gamma_d \left(1 \wedge \frac{\sigma}{A_{\max}} \right) \left(\frac{n_1 r}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq n_1, 1 \leq j \leq r \right\}, \quad (25)$$

and

$$\mathcal{A} = \left\{ \mathbf{A} = (a_{ij}) \in \mathbb{R}^{r \times n_2} : a_{ij} \in \left\{ 0, A_{\max}, \gamma_a(\sigma \wedge A_{\max}) \left(\frac{k}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq r, 1 \leq j \leq n_2 \text{ \& } \|\mathbf{A}\|_0 \leq k \right\}. \quad (26)$$

Clearly \mathcal{X} as defined in (24) is a finite class of matrices which admits a factorization as in Section II-A which implies that $\mathcal{X} \subset \mathcal{X}(r, k, A_{\max})$. We consider the lower bounds involving the non-sparse term, \mathbf{D} and the sparse factor \mathbf{A} separately and then combine those results to get an overall lower bound on the minimax risk $\mathcal{R}_{\mathcal{X}(r, k, A_{\max})}^*$.

Let us first establish the lower bound obtained by using the sparse factor \mathbf{A} . For this let us define a finite class $\mathcal{X}_A \subseteq \mathcal{X}$ as

$$\mathcal{X}_A \triangleq \left\{ \mathbf{X} = \mathbf{D}\mathbf{A} : \mathbf{D} = \left(\mathbf{I}_r \mid \cdots \mid \mathbf{I}_r \mid \mathbf{0}_A \right)^T \in \mathcal{D}, \mathbf{A} \in \bar{\mathcal{A}} \right\}, \quad (27)$$

where, \mathbf{I}_r denotes the $r \times r$ identity matrix, $\mathbf{0}_A$ is the $r \times (n_1 - \lfloor \frac{n_1}{r} \rfloor r)$ zero matrix and $\bar{\mathcal{A}} \subseteq \mathcal{A}$ is defined as

$$\bar{\mathcal{A}} = \left\{ \mathbf{A} = (a_{ij}) \in \mathbb{R}^{r \times n_2} : a_{ij} \in \left\{ 0, \gamma_a(\sigma \wedge A_{\max}) \left(\frac{k}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq r, 1 \leq j \leq n_2 \text{ \& } \|\mathbf{A}\|_0 \leq k \right\}. \quad (28)$$

The definition in (27) implies that the elements of \mathcal{X}_A are in the form of a block matrix $\left(\mathbf{A} \mid \cdots \mid \mathbf{A} \mid \mathbf{0}_A \right)$, $\forall \mathbf{A} \in \bar{\mathcal{A}}$. Overall there are $\lfloor n_1/r \rfloor$ blocks of \mathbf{A} and the rest is a zero matrix of the appropriate dimension. Since the entries of \mathbf{A} can take only one of two values 0 or $\left[\gamma_a(\sigma \wedge A_{\max}) \sqrt{k/m} \right]$ and since there are at most k non-zero elements (due to the sparsity constraint), the Varshamov-Gilbert bound (cf. Lemma 2.9 in [21]) guarantees the existence of a subset $\mathcal{X}_A^0 \subset \mathcal{X}_A$ with cardinality $\text{Card}(\mathcal{X}_A^0) \geq 2^{k/8} + 1$, containing the $n_1 \times n_2$ zero matrix $\mathbf{0}$, such that for any 2 distinct elements $\mathbf{X}_1, \mathbf{X}_2 \in \mathcal{X}_A^0$ we have,

$$\begin{aligned} \|\mathbf{X}_1 - \mathbf{X}_2\|_F^2 &\geq \left(\frac{k}{8} \right) \left\lfloor \frac{n_1}{r} \right\rfloor \gamma_a^2(\sigma \wedge A_{\max})^2 \left(\frac{k}{m} \right) \\ &\geq \frac{\gamma_a^2}{16} (\sigma \wedge A_{\max})^2 \left(\frac{n_1 k}{r} \right) \left(\frac{k}{m} \right), \end{aligned} \quad (29)$$

where the last inequality comes from the fact that $\lfloor x \rfloor \geq x/2 \forall x \geq 1$.

Now, since the observations are corrupted by independent zero-mean additive Gaussian noise $\xi_{i,j} \sim \mathcal{N}(0, \sigma^2)$, the joint pdf of the set of $|\mathcal{S}|$ measurements $\mathbf{Y}_{\mathcal{S}}$ (conditioned on $\mathbf{X}_{\mathcal{S}}^*$) is a multivariate Gaussian density of dimension $|\mathcal{S}|$ whose mean is the collection of true matrix entries at the sample locations, and covariance matrix $\sigma^2 \mathbf{I}_{|\mathcal{S}|}$, where $\mathbf{I}_{|\mathcal{S}|}$ is the $|\mathcal{S}| \times |\mathcal{S}|$ identity matrix. Hence,

$$p_{\mathbf{X}_{\mathcal{S}}^*}(\mathbf{Y}_{\mathcal{S}}) = \frac{1}{(2\pi\sigma^2)^{|\mathcal{S}|/2}} \exp \left(-\frac{1}{2\sigma^2} \sum_{(i,j) \in \mathcal{S}} (Y_{i,j} - X_{i,j}^*)^2 \right). \quad (30)$$

Using the expression for the joint pdf in (30), for any $\mathbf{X} \in \mathcal{X}_A^0$ and a sampling set \mathcal{S} , the KL divergence of \mathbb{P}_0

from $\mathbb{P}_{\mathbf{X}}$ satisfies

$$\begin{aligned} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) &= \mathbb{E}_{\mathbf{X}} \left[\log \frac{p_{\mathbf{X}_S}(\mathbf{Y}_S)}{p_0(\mathbf{Y}_S)} \right] \\ &= \frac{m}{n_1 n_2} \left(\frac{1}{2\sigma^2} \right) \sum_{i,j} |X_{i,j}|^2 \end{aligned} \quad (31)$$

$$\leq \frac{m}{2\sigma^2} \gamma_a^2 (\sigma \wedge A_{\max})^2 \left(\frac{k}{m} \right) \quad (32)$$

where (31) is obtained by conditioning¹ the expectation w.r.t the sampling set \mathcal{S} . From (32) we see that

$$\frac{1}{\text{Card}(\mathcal{X}_A^0 - 1)} \sum_{\mathbf{X} \in \mathcal{X}_A^0} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) \leq \alpha \log(\text{Card}(\mathcal{X}_A^0) - 1) \quad (33)$$

is satisfied for any $\alpha > 0$ by choosing $0 < \gamma_a < \frac{\sigma \sqrt{\alpha \log 2}}{2(\sigma \wedge A_{\max})}$. Equations (29) and (33) imply we can apply Theorem A.1 (where the Frobenius error has been as used the semi-distance function) to yield

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}_A} \mathbb{P}_{\mathbf{X}^*} \left(\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} \geq \frac{\gamma_a^2}{32} (\sigma \wedge A_{\max})^2 \left(\frac{k}{m} \right) \left(\frac{k}{m n_2} \right) \right) \geq \beta \quad (34)$$

for some absolute constant $\beta \in (0, 1)$.

We now consider the non-sparse factor \mathbf{D} to construct a testing set and establish lower bounds similar to the previous case. In order to obtain tight lower bounds we shall examine this scenario in two different sparsity regimes (as mentioned in Section III). First let us define a finite class of matrices $\mathcal{X}_D \subseteq \mathcal{X}$ as

$$\mathcal{X}_D \triangleq \left\{ \mathbf{X} = \mathbf{D}\mathbf{A} : \mathbf{D} \in \bar{\mathcal{D}}, \mathbf{A} = A_{\max} \left(\mathbf{I}_r \mid \cdots \mid \mathbf{I}_r \mid \Psi_D \right) \in \mathcal{A}, \right\}, \quad (35)$$

where, \mathbf{I}_r denotes the $r \times r$ identity matrix, Ψ_D takes different forms in the two sparsity regimes:

- For $k < n_2$: Ψ_D is the $r \times (n_2 - r \lfloor k/r \rfloor)$ zero matrix
- For $k \geq n_2$: $\Psi_D = \begin{pmatrix} \mathbf{I}_D \\ \mathbf{0}_D \end{pmatrix}$, where, \mathbf{I}_D is the identity matrix of dimension $(n_2 - r \lfloor n_2/r \rfloor)$ and $\mathbf{0}_D$ is a zero matrix of dimension $(r - n_2 + r \lfloor n_2/r \rfloor) \times (n_2 - r \lfloor n_2/r \rfloor)$

and $\bar{\mathcal{D}} \subseteq \mathcal{D}$ is defined as follows

$$\bar{\mathcal{D}} = \left\{ \mathbf{D} = (d_{ij}) \in \mathbb{R}^{n_1 \times r} : d_{ij} \in \left\{ 0, \gamma_d \left(1 \wedge \frac{\sigma}{A_{\max}} \right) \left(\frac{n_1 r}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq n_1, 1 \leq j \leq r \right\}. \quad (36)$$

The definition in (35) is similar to that we used to construct \mathcal{X}_A and hence it results in a block matrix structure for the elements in \mathcal{X}_D . We note here that there are $n_1 r$ elements in each block \mathbf{D} , where each entry can be either 0 or $\left[\gamma_d \left(1 \wedge \frac{\sigma}{A_{\max}} \right) \left(\frac{n_1 r}{m} \right)^{1/2} \right]$. Hence the Varshamov-Gilbert bound (cf. Lemma 2.9 in [21]) guarantees the existence of a subset $\mathcal{X}_D^0 \subset \mathcal{X}_D$ with cardinality $\text{Card}(\mathcal{X}_D^0) \geq 2^{n_1 r/8} + 1$, containing the $n_1 \times n_2$ zero matrix $\mathbf{0}$, such that for any 2 distinct elements $\mathbf{X}_1, \mathbf{X}_2 \in \mathcal{X}_D^0$ we have

$$\|\mathbf{X}_1 - \mathbf{X}_2\|_F^2 \geq \left(\frac{n_1 r}{8} \right) \gamma_d^2 (\sigma \wedge A_{\max})^2 \left(\frac{n_1 r}{m} \right) \delta(k, n_2), \quad (37)$$

¹Here, both the observations \mathbf{Y}_S , and the sampling set \mathcal{S} are random quantities. Thus by conditioning w.r.t to \mathcal{S} , we get $\mathbb{E}_{\mathbf{X}} \triangleq \mathbb{E}_{\mathbf{Y}_S \sim p_{\mathbf{X}_S}} [\cdot] = \mathbb{E}_{\mathcal{S}} [\mathbb{E}_{\mathbf{X}_S | \mathcal{S}} [\cdot]]$. Since \mathcal{S} is generated according to the independent Bernoulli($m/n_1 n_2$) model, $\mathbb{E}_{\mathcal{S}} [\cdot]$ yields the constant term $\frac{m}{n_1 n_2}$. We shall use such conditioning techniques in subsequent proofs as well.

where $\delta(k, n_2)$ is given by

$$\delta(k, n_2) = \begin{cases} \lfloor k/r \rfloor & \text{for } k < n_2 \\ (n_2/r) & \text{for } k \geq n_2 \end{cases}. \quad (38)$$

Using (38) in (37) we get

$$\|\mathbf{X}_1 - \mathbf{X}_2\|_F^2 \geq \frac{\gamma_d^2 n_1 n_2}{16} (\sigma \wedge A_{\max})^2 \left(\frac{n_1 r}{m}\right) \Delta(k, n_2), \quad (39)$$

where $\Delta(k, n_2)$ is as in (8).

Using the expression for the joint pdf in (30), for any $\mathbf{X} \in \mathcal{X}_D^0$ and a sampling set \mathcal{S} , the KL divergence of \mathbb{P}_0 from $\mathbb{P}_{\mathbf{X}}$ satisfies

$$\begin{aligned} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) &= \mathbb{E}_{\mathbf{X}} \left[\log \frac{p_{\mathbf{X}_{\mathcal{S}}}(\mathbf{Y}_{\mathcal{S}})}{p_0(\mathbf{Y}_{\mathcal{S}})} \right] \\ &= \frac{m}{n_1 n_2} \left(\frac{1}{2\sigma^2} \right) \sum_{i,j} |X_{i,j}|^2 \\ &\leq \frac{m}{2\sigma^2} \gamma_d^2 (\sigma \wedge A_{\max})^2 \left(\frac{n_1 r}{m} \right) \end{aligned} \quad (40)$$

From (40) we see that

$$\frac{1}{\text{Card}(\mathcal{X}_D^0) - 1} \sum_{\mathbf{X} \in \mathcal{X}_D^0} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) \leq \alpha' \log(\text{Card}(\mathcal{X}_D^0) - 1) \quad (41)$$

is satisfied for any $\alpha' > 0$ by choosing $0 < \gamma_d < \frac{\sigma \sqrt{\alpha' \log 2}}{2(\sigma \wedge A_{\max})}$. Equations (39) and (41) imply we can apply Theorem A.1 (where the Frobenius error has been as used the semi-distance function) to yield

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}_D} \mathbb{P}_{\mathbf{X}^*} \left(\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} \geq \frac{\gamma_d^2}{32} (\sigma \wedge A_{\max})^2 \left(\frac{n_1 r}{m} \right) \Delta(k, n_2) \right) \geq \beta', \quad (42)$$

for some absolute constant $\beta' \in (0, 1)$. Inequalities (34) and (42) imply the result,

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}(r, k, A_{\max})} \mathbb{P}_{\mathbf{X}^*} \left(\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} \geq C (\sigma \wedge A_{\max})^2 \left[\left(\frac{n_1 r}{m} \right) \Delta(k, n_2) + \left(\frac{k}{m} \right) \left(\frac{k}{r n_2} \right) \right] \right) \geq (\beta' \vee \beta) \quad (43)$$

where $C = \frac{(\gamma_d \wedge \gamma_a)^2}{64}$, is a suitable value for the leading constant, and we have $(\beta' \vee \beta) \in (0, 1)$. In order to obtain this result for the entire class $\mathcal{X}(r, k, A_{\max})$, we have used the fact that solving the problem over a larger (and possibly uncountable) class of matrices is at least as tough as solving the same problem over a smaller (and possibly finite) subclass. Applying Markov's inequality to (43) directly yields the result of Theorem III.1, completing the proof. ■

B. Proof of Theorem III.2 (Sketch)

We follow a similar approach as in the previous proof, by establishing lower bounds involving the non-sparse term \mathbf{D} and the sparse factor \mathbf{A} separately and then combining them to get the final result. The arguments used here will be analogous to those in Appendix A but for the definitions of the finite classes defined and the way the KL divergences are handled. To avoid repetition, we present a sketch highlighting how this differs from the one in Appendix A.

For a finite class of matrices $\mathcal{X} \subset \mathbb{R}^{n_1 \times n_2}$ as constructed in (24), the factor classes \mathcal{D} and \mathcal{A} are defined by

$$\mathcal{D} = \left\{ \mathbf{D} = (d_{ij}) \in \mathbb{R}^{n_1 \times r} : d_{ij} \in \left\{ 0, 1, \gamma_d (1 \wedge (\tau A_{\max})^{-1}) \left(\frac{n_1 r}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq n_1, 1 \leq j \leq r \right\} \quad (44)$$

and

$$\mathcal{A} = \left\{ \mathbf{A} = (a_{ij}) \in \mathbb{R}^{r \times n_2} : a_{ij} \in \left\{ 0, A_{\max}, \gamma_a (\tau^{-1} \wedge A_{\max}) \left(\frac{k}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq r, 1 \leq j \leq n_2, \|\mathbf{A}\|_0 \leq k \right\}. \quad (45)$$

Now, since the observations are corrupted with independent zero-mean additive Laplace noise (denoted $\xi_{i,j} \sim \text{Laplace}(0, \tau)$), $\forall (i, j) \in \mathcal{S}$ where $\tau > 0$, the joint pdf of the set of $|\mathcal{S}|$ measurements $\mathbf{Y}_{\mathcal{S}}$ (conditioned on $\mathbf{X}_{\mathcal{S}}^*$) is a multivariate Laplace distribution given by

$$p_{\mathbf{X}_{\mathcal{S}}^*}(\mathbf{Y}_{\mathcal{S}}) = \left(\frac{\tau}{2} \right)^{|\mathcal{S}|} \exp \left(-\tau \sum_{(i,j) \in \mathcal{S}} |Y_{i,j} - X_{i,j}^*| \right) \quad (46)$$

Using the expression for the joint pdf in (46), for any $\mathbf{X} \in \mathcal{X}_A^0$ (or \mathcal{X}_D^0) and a known sampling set \mathcal{S} , the KL divergence of \mathbb{P}_0 from $\mathbb{P}_{\mathbf{X}}$ satisfies

$$\begin{aligned} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) &= \mathbb{E}_{\mathbf{X}} \left[\log \frac{p_{\mathbf{X}_{\mathcal{S}}}(\mathbf{Y}_{\mathcal{S}})}{p_0(\mathbf{Y}_{\mathcal{S}})} \right] \\ &= \sum_{i,j} K(\mathbb{P}_{X_{i,j}}, \mathbb{P}_{0_{i,j}}) \cdot \frac{m}{n_1 n_2} \end{aligned} \quad (47)$$

$$\leq \frac{\tau^2}{2} \sum_{i,j} |X_{i,j}|^2 \cdot \frac{m}{n_1 n_2}, \quad (48)$$

where we get (47) by conditioning the expectation w.r.t the sampling set \mathcal{S} , and (48) is obtained using the following lemma.

Lemma A.1. *For a Laplace distribution with parameter τ centered at x , denoted $\mathbb{P}_x \sim \text{Laplace}(x, \tau)$ where $x \in \mathbb{R}$ and $\tau > 0$, the KL divergence of \mathbb{P}_x from \mathbb{P}_y (for any $y \in \mathbb{R}$) satisfies*

$$K(\mathbb{P}_x, \mathbb{P}_y) \leq \frac{\tau^2}{2} (x - y)^2.$$

Proof: For \mathbb{P}_x and \mathbb{P}_y as defined above we have, by (relatively) straightforward calculation

$$\begin{aligned} K(\mathbb{P}_x, \mathbb{P}_y) &= \mathbb{E}_x \left[\log \frac{p_x(z)}{p_y(z)} \right] \\ &= \tau |x - y| - (1 - e^{-\tau |x - y|}). \end{aligned} \quad (49)$$

Using a series expansion of the exponent in (49) we have

$$\begin{aligned} e^{-\tau |x - y|} &= 1 - \tau |x - y| + \frac{(\tau |x - y|)^2}{2!} - \frac{(\tau |x - y|)^3}{3!} + \dots \\ &\leq 1 - \tau |x - y| + \frac{(\tau |x - y|)^2}{2!}. \end{aligned} \quad (50)$$

Rearranging the terms in (50) yields the result. ■

With the upper bound on the KL divergence² established in (48) and the classes \mathcal{D} and \mathcal{A} defined as in Equations (44) and (45) respectively, the rest of the arguments follow exactly as in the previous proof and hence are omitted here. ■

C. Proof of Theorem III.3

The Poisson observation model considered here assumes that all the entries of the underlying matrix \mathbf{X} are strictly non-zero. A straightforward observation that follows is that the sparse factor \mathbf{A} in the factorization cannot have any zero valued columns. Hence we have that $k \geq n_2$ be satisfied as a necessary (but not a sufficient) condition. We shall use similar techniques as in the previous sections to derive the result for this model. However we need to be careful while constructing the sample class of matrices as we need to ensure that all the entries of the members should be strictly bounded away from zero (and in fact $\geq X_{\min}$).

For some $\delta \in (0, 1)$, let us define a matrix $\mathbf{D}_\delta \in \mathbb{R}^{n_1 \times r}$ such that $(\mathbf{D}_\delta)_{ij} = \delta, \forall 1 \leq i \leq n_1, 1 \leq j \leq r$. Using this we define a class of matrices $\mathcal{X} \subset \mathcal{X}'(r, k, A_{\max})$ as

$$\mathcal{X} \triangleq \{\mathbf{X} = (\mathbf{D}_\delta + \mathbf{D})\mathbf{A} : \mathbf{D} \in \mathcal{D}, \mathbf{A} \in \mathcal{A}\}, \quad (51)$$

where the factor classes $\mathcal{D} \subset \mathbb{R}^{n_1 \times r}$ and $\mathcal{A} \subset \mathbb{R}^{r \times n_2}$ are defined as follows for some $\gamma_d, \gamma_a \leq 1$ (which we shall qualify later)

$$\mathcal{D} = \left\{ \mathbf{D} = (d_{ij}) \in \mathbb{R}^{n_1 \times r} : d_{ij} \in \left\{ 0, 1 - \delta, \gamma_d \left(1 \wedge \frac{X_{\min} + \sqrt{X_{\min}}}{A_{\max}} \right) \left(\frac{n_1 r}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq n_1, 1 \leq j \leq r \right\} \quad (52)$$

and

$$\mathcal{A} = \left\{ \mathbf{A} = (a_{ij}) \in \mathbb{R}^{r \times n_2} : a_{ij} \in \left\{ 0, A_{\max}, \gamma_a \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right] \left(\frac{k - n_2}{m} \right)^{1/2} \right\}, \right. \\ \left. \forall 1 \leq i \leq r, 1 \leq j \leq n_2 \text{ \& } \|\mathbf{A}\|_0 \leq k \right\}. \quad (53)$$

Clearly \mathcal{X} as defined in (51) is a finite class of matrices which admits a factorization as in Section II-A which implies that $\mathcal{X} \subset \mathcal{X}'(r, k, A_{\max})$. As before, we consider the lower bounds involving the non-sparse term, \mathbf{D} and the sparse factor \mathbf{A} separately and then combine those results to get an overall lower bound on the minimax risk $\mathcal{R}_{\mathcal{X}'(r, k, A_{\max})}^*$.

²The KL divergence of \mathbb{P}_0 from $\mathbb{P}_{\mathbf{X}}$ also satisfies $K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_0) = \sum_{(i,j) \in \mathcal{S}} \left[\tau |X_{i,j} - 0| - (1 - e^{-\tau |X_{i,j} - 0|}) \right] \leq \tau \sum_{(i,j) \in \mathcal{S}} |X_{i,j}|$, where the final inequality follows from the fact that $(1 - e^{-\tau |X_{i,j}|}) \geq 0$. Using this linear upper bound on the KL divergence, we may obtain lower bounds on the minimax risk of the form

$$\mathcal{R}_{\mathcal{X}'(r, k, A_{\max})}^* \geq c (\tau^{-1} \wedge A_{\max})^2 \left\{ \left(\frac{n_1 r}{m} \right)^2 \Delta(k, n_2) + \left(\frac{k}{m} \right)^2 \left(\frac{k}{r n_2} \right) \right\}.$$

When the expected number of samples is low i.e. when $m \leq (n_1 r \wedge k)$, the above bound is stronger than the result in Theorem III.2. Using similar analysis it is easy to strengthen the lower bounds of the minimax risk in other sampling regimes.

Let us first establish the lower bound obtained by using the sparse factor \mathbf{A} . For this let us define a finite class $\mathcal{X}_A \subseteq \mathcal{X}$ as

$$\mathcal{X}_A \triangleq \left\{ \mathbf{X} = (\mathbf{D}_\delta + \mathbf{D}_I) \mathbf{A} : \mathbf{D}_I = (1 - \delta) \left(\mathbf{I}_r \mid \cdots \mid \mathbf{I}_r \mid \Psi_A \right)^T \in \mathcal{D}, \mathbf{A} \in \bar{\mathcal{A}} \right\}, \quad (54)$$

and \mathbf{I}_r denotes the $r \times r$ identity matrix, and Ψ_A is a $r \times (n_1 - \lfloor \frac{n_1}{r} \rfloor r)$ matrix given by $\Psi_A = \begin{pmatrix} \mathbf{I}_A \\ \mathbf{0}_A \end{pmatrix}$, where \mathbf{I}_A is the identity matrix of dimension $(n_1 - r \lfloor n_1/r \rfloor)$ and $\mathbf{0}_A$ is the $(r - n_1 + r \lfloor n_1/r \rfloor) \times (n_1 - r \lfloor n_1/r \rfloor)$ zero matrix and the class $\bar{\mathcal{A}} \subseteq \mathcal{A}$ consists of a special class of matrices such that $\forall \mathbf{A} \in \bar{\mathcal{A}}$ we have,

$$(\mathbf{A})_{ij} = \begin{cases} A_{\max} & \text{for } i = 1, 1 \leq j \leq n_2 \\ (\mathbf{A})_{ij} \in \left\{ 0, \gamma_a \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right] \left(\frac{k-n_2}{m} \right)^{1/2} \right\} & \text{for } 2 \leq i \leq r, 1 \leq j \leq n_2 \end{cases}. \quad (55)$$

The definition in (54) implies that the entries of any matrix $\mathbf{X} \in \mathcal{X}_A$ satisfy $\min_{(i,j) \in [n_1] \times [n_2]} X_{i,j} \geq X_{\min}$ provided we choose δ such that $\delta A_{\max} \geq X_{\min}$. Now let us consider the Frobenius distance between any two distinct elements $\mathbf{X}_1, \mathbf{X}_2 \in \mathcal{X}_A$,

$$\begin{aligned} \|\mathbf{X}_1 - \mathbf{X}_2\|_F^2 &= \|(\mathbf{D}_I + \mathbf{D}_\delta) \mathbf{A}_1 - (\mathbf{D}_I + \mathbf{D}_\delta) \mathbf{A}_2\|_F^2 \\ &= \|\mathbf{D}_I (\mathbf{A}_1 - \mathbf{A}_2) + \mathbf{D}_\delta (\mathbf{A}_1 - \mathbf{A}_2)\|_F^2 \\ &= \|\mathbf{D}_I (\mathbf{A}_1 - \mathbf{A}_2)\|_F^2 + \underbrace{\|\mathbf{D}_\delta (\mathbf{A}_1 - \mathbf{A}_2)\|_F^2}_{\geq 0} + \underbrace{2 \cdot \text{trace}(\mathbf{D}_\delta (\mathbf{A}_1 - \mathbf{A}_2) (\mathbf{A}_1 - \mathbf{A}_2)^T \mathbf{D}_I^T)}_{\geq 0} \\ &\geq \|\mathbf{D}_I \mathbf{A}_1 - \mathbf{D}_I \mathbf{A}_2\|_F^2. \end{aligned} \quad (56)$$

The definition of \mathbf{A} as in (55) and the sparsity constraint on the entries imply that the number of degrees of freedom in the sparse factor is restricted to $(k - n_2)$. The Varshamov-Gilbert bound (cf. Lemma 2.9 in [21]) can be easily applied to the set of matrices of the form $\tilde{\mathbf{X}} = \mathbf{D}_I \mathbf{A}$ where $\mathbf{A} \in \bar{\mathcal{A}}$, and this when coupled with (56) guarantees the existence of a subset $\mathcal{X}_A^0 \subset \mathcal{X}_A$ with cardinality $\text{Card}(\mathcal{X}_A^0) \geq 2^{(k-n_2)/8} + 1$, such that for any two distinct elements $\mathbf{X}_1, \mathbf{X}_2 \in \mathcal{X}_A^0$ we have,

$$\begin{aligned} \|\mathbf{X}_1 - \mathbf{X}_2\|_F^2 &\geq (1 - \delta)^2 \left(\frac{k - n_2}{8} \right) \left\lfloor \frac{n_1}{r} \right\rfloor \gamma_a^2 \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right]^2 \left(\frac{k - n_2}{m} \right) \\ &\geq (1 - \delta)^2 \frac{\gamma_a^2}{16} \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right]^2 \left(\frac{n_1(k - n_2)}{r} \right) \left(\frac{k - n_2}{m} \right). \end{aligned} \quad (57)$$

The joint pmf of the set of $|\mathcal{S}|$ observations (conditioned on the true underlying matrix) can be conveniently written as a product of Poisson pmfs using the independence criterion as,

$$p_{\mathbf{X}_S^*}(\mathbf{Y}_S) = \prod_{(i,j) \in \mathcal{S}} \frac{(X_{i,j}^*)^{Y_{i,j}} e^{-X_{i,j}^*}}{(Y_{i,j})!}. \quad (58)$$

For any $\mathbf{X} \in \mathcal{X}_A^0$ the KL divergence of $\mathbb{P}_{\mathbf{X}'}$ from $\mathbb{P}_{\mathbf{X}}$ where \mathbf{X}' is the reference matrix (which corresponds to

the zero matrix in the other models) is obtained by using an intermediate result from [23] giving

$$\begin{aligned} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_{\mathbf{X}'}) &= \mathbb{E}_{\mathcal{S}} \left[\sum_{i,j} K(\mathbb{P}_{X_{i,j}}, \mathbb{P}_{X'_{i,j}}) \right] \\ &= \frac{m}{n_1 n_2} \sum_{i,j} \left\{ X_{i,j} \log \left(\frac{X_{i,j}}{X'_{i,j}} \right) - X_{i,j} + X'_{i,j} \right\}. \end{aligned}$$

Using the inequality $\log t \leq (t - 1)$, we can bound the KL divergence as

$$\begin{aligned} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_{\mathbf{X}'}) &\leq \frac{m}{n_1 n_2} \sum_{i,j} \left\{ X_{i,j} \left(\frac{X_{i,j} - X'_{i,j}}{X'_{i,j}} \right) - X_{i,j} + X'_{i,j} \right\} \\ &\leq \frac{m}{n_1 n_2} \sum_{i,j} \frac{(X_{i,j} - X'_{i,j})^2}{X'_{i,j}} \\ &\leq m \frac{(X_{i,j} - X_{\min})^2}{X_{\min}}. \end{aligned} \quad (59)$$

From (59) and the bound on the entries of the elements in \mathcal{X}_A^0 we see that

$$\frac{1}{\text{Card}(\mathcal{X}_A^0 - 1)} \sum_{\mathbf{X} \in \mathcal{X}_A^0} K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_{\mathbf{X}'}) \leq \alpha \log(\text{Card}(\mathcal{X}_A^0) - 1) \quad (60)$$

is satisfied for any $\alpha > 0$ by choosing $\max \left\{ \left(\frac{\sqrt{8X_{\min}} - \sqrt{\alpha \log 2}}{\sqrt{8(\sqrt{X_{\min}} + 1)}} \right), 0 \right\} < \gamma_a < \left(\frac{\sqrt{8X_{\min}} + \sqrt{\alpha \log 2}}{\sqrt{8(\sqrt{X_{\min}} + 1)}} \right)$. Equations (57) and (60) imply we can apply Theorem A.1 (where the Frobenius error has been as used the semi-distance function) to yield

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}_A} \mathbb{P}_{\mathbf{X}^*} \left(\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} \geq (1 - \delta)^2 \frac{\gamma_a^2}{32} \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right]^2 \left(\frac{k - n_2}{m} \right) \left(\frac{k - n_2}{r n_2} \right) \right) \geq \beta \quad (61)$$

for some absolute constant $\beta \in (0, 1)$. Establishing lower bounds using the dictionary term \mathbf{D} follows a similar and more direct route than the previous case. For this let us first define the finite class, $\mathcal{X}_D \subseteq \mathcal{X}$ as

$$\mathcal{X}_D \triangleq \left\{ \mathbf{X} = (\mathbf{D}_\delta + \mathbf{D})\mathbf{A} : \mathbf{D} \in \bar{\mathcal{D}}, \mathbf{A} = A_{\max} \left(\mathbf{I}_r \mid \cdots \mid \mathbf{I}_r \mid \Psi_D \right) \in \mathcal{A}, \right\}, \quad (62)$$

where, \mathbf{I}_r denotes the $r \times r$ identity matrix, Ψ_D is a $r \times (n_2 - r \lfloor k/r \rfloor)$ matrix given by $\Psi_D = \begin{pmatrix} \mathbf{I}_D \\ \mathbf{0}_D \end{pmatrix}$ and, \mathbf{I}_D is the identity matrix of dimension $(n_2 - r \lfloor n_2/r \rfloor)$ and $\mathbf{0}_D$ is the $(r - n_2 + r \lfloor n_2/r \rfloor) \times (n_2 - r \lfloor n_2/r \rfloor)$ zero matrix, and $\bar{\mathcal{D}} \subseteq \mathcal{D}$ is defined as

$$\bar{\mathcal{D}} = \left\{ \mathbf{D} = (d_{ij}) \in \mathbb{R}^{n_1 \times r} : d_{ij} \in \left\{ 0, \gamma_d \left(1 \wedge \frac{X_{\min} + \sqrt{X_{\min}}}{A_{\max}} \right) \left(\frac{n_1 r}{m} \right)^{1/2} \right\}, \forall 1 \leq i \leq n_1, 1 \leq j \leq r \right\}. \quad (63)$$

The definition of \mathcal{X}_D as in (62) enables us to invoke the Varshamov-Gilbert bound (cf. Lemma 2.9 in [21]) and the KL divergences are bounded using similar arguments as in the previous case. To avoid repetition, we omit the details here and state the result for the lower bound on the minimax risk using the dictionary term

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X}^* \in \mathcal{X}_A} \mathbb{P}_{\mathbf{X}^*} \left(\frac{\|\hat{\mathbf{X}} - \mathbf{X}^*\|_F^2}{n_1 n_2} \geq (1 - \delta)^2 \frac{\gamma_d^2}{32} \left[(X_{\min} + \sqrt{X_{\min}}) \wedge A_{\max} \right]^2 \left(\frac{n_1 r}{m} \right) \right) \geq \beta' \quad (64)$$

for some absolute constant $\beta' \in (0, 1)$. Using (61) and (64) and by applying Markov's inequality, we directly get the result presented in Theorem III.3 thus completing the proof. \blacksquare

D. Proof of Theorem III.4 (Sketch)

For any \mathbf{X} , $\mathbf{X}^* \in \mathcal{X}(r, k, A_{\max})$ using the pdf model described in (16), it is straightforward to show that the scalar KL divergence is given by

$$K(\mathbb{P}_{X_{i,j}^*}, \mathbb{P}_{X_{i,j}}) = F(X_{i,j}^*) \log \left(\frac{F(X_{i,j}^*)}{F(X_{i,j})} \right) + (1 - F(X_{i,j}^*)) \log \left(\frac{1 - F(X_{i,j}^*)}{1 - F(X_{i,j})} \right) \quad (65)$$

for any $(i, j) \in \mathcal{S}$. We directly use an intermediate result from [20] to invoke a quadratic upper bound for the KL divergence term,

$$K(\mathbb{P}_{X_{i,j}^*}, \mathbb{P}_{X_{i,j}}) \leq \frac{1}{2} c_{F, rA_{\max}} (X_{i,j}^* - X_{i,j})^2, \quad (66)$$

where $c_{F, rA_{\max}}$ is defined in (17). Such an upper bound in terms of the underlying matrix entries can be attained by following a procedure illustrated in [26], where one first establishes quadratic bounds on the KL divergence in terms of the Bernoulli parameters, and then subsequently establishes a bound on the squared difference between Bernoulli parameters in terms of the squared difference of the underlying matrix elements.

Using (66) and the independence of $\mathbf{Y}_{\mathcal{S}}$ (conditioned on $\mathbf{X}_{\mathcal{S}}$), for any $\mathbf{X} \in \mathcal{X}(r, k, A_{\max})$ and the $n_1 \times n_2$ zero-matrix $\mathbf{0}$, the KL divergence of $\mathbb{P}_{\mathbf{0}}$ from $\mathbb{P}_{\mathbf{X}}$ obeys the bound

$$K(\mathbb{P}_{\mathbf{X}}, \mathbb{P}_{\mathbf{0}}) \leq \frac{m}{2n_1 n_2} c_{F, rA_{\max}} \sum_{i,j} X_{i,j}^2. \quad (67)$$

With a quadratic upper bound on the KL divergence in terms of the underlying matrix entries (67), the rest of the arguments follow in an analogous manner as in the Gaussian case, and hence are omitted. ■

REFERENCES

- [1] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Foundations of Computational mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [2] E. J. Candès and T. Tao, “The power of convex relaxation: Near-optimal matrix completion,” *IEEE Trans. Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [3] R. Keshavan, A. Montanari, and S. Oh, “Matrix completion from a few entries,” *IEEE Trans. Information Theory*, vol. 56, no. 6, pp. 2980–2998, 2010.
- [4] B. Recht, “A simpler approach to matrix completion,” *The Journal of Machine Learning Research*, vol. 12, pp. 3413–3430, 2011.
- [5] D. Gross, “Recovering low-rank matrices from few coefficients in any basis,” *IEEE Trans. Information Theory*, vol. 57, no. 3, pp. 1548–1566, 2011.
- [6] R. Keshavan, A. Montanari, and S. Oh, “Matrix completion from noisy entries,” in *Advances in Neural Information Processing Systems*, 2009, pp. 952–960.
- [7] E. J. Candès and Y. Plan, “Matrix completion with noise,” *Proceedings of the IEEE*, vol. 98, no. 6, pp. 925–936, 2010.
- [8] V. Koltchinskii, K. Lounici, and A. B. Tsybakov, “Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion,” *The Annals of Statistics*, vol. 39, no. 5, pp. 2302–2329, 2011.
- [9] A. Rohde and A. B. Tsybakov, “Estimation of high-dimensional low-rank matrices,” *The Annals of Statistics*, vol. 39, no. 2, pp. 887–930, 2011.
- [10] T. T. Cai and W. Zhou, “Matrix completion via max-norm constrained optimization,” *arXiv preprint arXiv:1303.0341*, 2013.
- [11] O. Klopp, “Noisy low-rank matrix completion with general sampling distribution,” *Bernoulli*, vol. 20, no. 1, pp. 282–303, 2014.
- [12] J. Lafond, “Low rank matrix completion with exponential family noise,” *arXiv preprint arXiv:1502.06919*, 2015.
- [13] M. A. Davenport, Y. Plan, E. van den Berg, and M. Wootters, “1-bit matrix completion,” *Information and Inference*, vol. 3, no. 3, pp. 189–223, 2014.

- [14] Y. Plan, R. Vershynin, and E. Yudovina, “High-dimensional estimation with geometric constraints,” *arXiv preprint arXiv:1404.3749*, 2014.
- [15] A. Soni and J. Haupt, “Estimation error guarantees for poisson denoising with sparse and structured dictionary models,” in *IEEE Intl Symposium on Information Theory*. IEEE, 2014, pp. 2002–2006.
- [16] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?,” *Journal of the ACM (JACM)*, vol. 58, no. 3, pp. 11, 2011.
- [17] D. Hsu, S. M. Kakade, and T. Zhang, “Robust matrix decomposition with sparse corruptions,” *IEEE Trans. Information Theory*, vol. 57, no. 11, pp. 7221–7234, 2011.
- [18] H. Xu, C. Caramanis, and S. Sanghavi, “Robust pca via outlier pursuit,” *IEEE Trans. Information Theory*, vol. 58, no. 5, pp. 3047–3064, 2012.
- [19] Y. Chen, A. Jalali, S. Sanghavi, and C. Caramanis, “Low-rank matrix recovery from errors and erasures,” *IEEE Trans. Information Theory*, vol. 59, no. 7, pp. 4324–4337, 2013.
- [20] A. Soni, S. Jain, J. Haupt, and S. Gonella, “Noisy matrix completion under sparse factor models,” *arXiv preprint arXiv:1411.0282*, 2014.
- [21] A. B. Tsybakov, *Introduction to nonparametric estimation*, Springer, 2008.
- [22] O. Klopp, K. Lounici, and A. B. Tsybakov, “Robust matrix completion,” *arXiv preprint arXiv:1412.8132*, 2014.
- [23] M. Raginsky, R. M. Willett, Z. T. Harmany, and R. F. Marcia, “Compressed sensing performance bounds under poisson noise,” *IEEE Trans. Signal Processing*, vol. 58, no. 8, pp. 3990–4002, 2010.
- [24] E. D. Kolaczyk and R. D. Nowak, “Multiscale likelihood analysis and complexity penalized estimation,” *Annals of Statistics*, pp. 500–527, 2004.
- [25] X. Jiang, G. Raskutti, and R. Willett, “Minimax optimal rates for poisson inverse problems with physical constraints,” *arXiv preprint arXiv:1403.6532*, 2014.
- [26] J. Haupt, N. Sidiropoulos, and G. Giannakis, “Sparse dictionary learning from 1-bit data,” in *IEEE Intl. Conference on Acoustics, Speech and Signal Processing*. IEEE, 2014, pp. 7664–7668.
- [27] A. Ribeiro and G. B. Giannakis, “Bandwidth-constrained distributed estimation for wireless sensor networks-part i: Gaussian case,” *IEEE Trans. Signal Processing*, vol. 54, no. 3, pp. 1131–1143, 2006.
- [28] Z. Q. Luo, “Universal decentralized estimation in a bandwidth constrained sensor network,” *IEEE Trans. Information Theory*, vol. 51, no. 6, pp. 2210–2219, 2005.
- [29] J. Lafond, O. Klopp, E. Moulines, and J. Salmon, “Probabilistic low-rank matrix completion on finite alphabets,” in *Advances in Neural Information Processing Systems*, 2014, pp. 1727–1735.
- [30] T. G. Kolda and B. W. Bader, “Tensor decompositions and applications,” *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.